

Benchmarking Summarization Methods for Scientific Abstracts: From Classical Models to LLMs

Fabio Baumgärtel ¹, Enrico Bono ^{1,2} , Paul Perco ^{1,3}  and Matthias Ley ^{1,2} *

¹ Delta4 GmbH, Vienna, Austria

² Division of Pediatric Nephrology and Gastroenterology, Department of Pediatrics and Adolescent Medicine, Comprehensive Center for Pediatrics, Medical University Vienna, Vienna, Austria

³ Department of Internal Medicine IV, Medical University Innsbruck, Innsbruck, Austria

* Correspondence: matthias.ley@delta4.ai

Abstract

A single paragraph of about 200 words maximum. For research articles, abstracts should give a pertinent overview of the work. We strongly encourage authors to use the following style of structured abstracts, but without headings: (1) Background: place the question addressed in a broad context and highlight the purpose of the study; (2) Methods: describe briefly the main methods or treatments applied; (3) Results: summarize the article's main findings; (4) Conclusions: indicate the main conclusions or interpretations. The abstract should be an objective representation of the article, it must not contain results which are not presented and substantiated in the main text and should not exaggerate the main conclusions.

Keywords: benchmarking; natural language processing; text summarization; large language models

1. Introduction

- foo [1]
 - bar
- why text summarization in biomedical domain is important
history of text summarization?

2. Materials and Methods

2.1. Gold-Standard Dataset

To establish a reliable benchmark for automatic summarization, we assembled a gold-standard dataset of 1,000 biomedical articles drawn from a diverse set of peer-reviewed journals hosted on *ScienceDirect* and *Cell Press*. These journals were selected because, in addition to their focus on molecular and biomedical sciences, they provide a standardized *Highlights* section [2,3]. This section provides concise bullet points that capture the main findings of each article. These served as the reference summaries in our evaluation, while the corresponding abstracts were used as input texts for the summarization.

Articles were collected systematically across a variety of journals to ensure coverage of different fields within molecular sciences such as drug discovery, genomics, proteomics, biotechnology, and biochemistry. We selected 50 articles from each of the 20 journals, bringing the dataset to 1,000 in total. The distribution of articles across journals is summarized in Table 1.

Received:

Revised:

Accepted:

Published:

Citation: . Title. *Int. J. Mol. Sci.* **2025**, *1*, 0. <https://doi.org/>

Copyright: © 2025 by the authors.

Submitted to *Int. J. Mol. Sci.* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Table 1. Overview of journals and number of articles included in the gold-standard dataset.

Publisher	Journal	Articles included
ScienceDirect	Drug Discovery Today	50
ScienceDirect	Journal of Molecular Biology	50
ScienceDirect	FEBS Letters	50
ScienceDirect	Journal of Biotechnology	50
ScienceDirect	Gene	50
ScienceDirect	Genomics	50
ScienceDirect	Journal of Proteomics	50
ScienceDirect	The International Journal of Biochemistry & Cell Biology	50
ScienceDirect	Cytokine	50
ScienceDirect	Developmental Cell	50
Cell	Cell	50
Cell	Cancer Cell	50
Cell	Cell Chemical Biology	50
Cell	Cell Genomics	50
Cell	Cell Host & Microbe	50
Cell	Cell Metabolism	50
Cell	Cell Reports	50
Cell	Cell Reports Medicine	50
Cell	Cell Stem Cell	50
Cell	Cell Systems	50

This setup provides standardized pairs of abstracts and reference summaries that can be directly used for evaluating automatic summarization methods.

2.2. Summarization Methods

We evaluated 50 summarization methods, ranging from simple frequency-based algorithms to state-of-the-art large language models (LLMs). By having this extensive coverage of methods, we were able to compare established techniques with the latest transformer-based models under identical conditions.

The models were grouped into five categories:

1. Traditional methods: As a foundation for comparison, we included two traditional extractive methods: a simple frequency-based approach and TextRank [4]. These methods provide a simple baseline to compare the more complex approaches with.
2. Encoder-Decoder models: We included a set of pre-trained encoder-decoder models, which are available through the HuggingFace library: BART (base and large) [5], T5 (base and large) [6], mT5 [7], and a variety of PEGASUS models [8]. These models are often applied for abstractive summarization and represent well-established neural systems within our benchmark.
3. General-purpose LLMs: We also evaluated a range of widely used large language models designed for broad application. This group includes models such as Gemma [9], Granite [10], LLaMA [11], Mistral [12], Phi [13,14], GPT [15,16], and Claude [17], which represent the current landscape of general-purpose systems.
4. Reasoning-oriented LLMs: We further included several models developed with a focus on advanced reasoning capabilities. This group includes models from the DeepSeek-R1 family [18], Qwen [19], more GPT models such as GPT-oss [20] and GPT-5 [21], Magistral [22], and some additional Claude models. Their design emphasizes multi-step problem solving and allowed us to explore whether reasoning affects summarization performance.
5. Specialized models: To assess whether domain adaptation improves summarization quality, we included MedLLaMA2 [23] (a medical adaptation of LLaMA-2) and LED

[24] (arXiv-tuned), which are trained on medical/biomedical data or on summarization tasks themselves.

The complete list of models included in each category is shown in Table 2.

Table 2. Overview of summarization methods/models evaluated in this study, organized by category.

Group	Methods/Models
Traditional methods	textrank; frequency
Encoder-Decoder models	facebook/bart-large-cnn; facebook/bart-base; google-t5/t5-base; google-t5/t5-large; cse-buethnlp/mT5_multilingual_XLSum; google/pegasus-xsum; google/pegasus-large; google/pegasus-cnn_dailymail
General-purpose LLMs	gemma3:1b; gemma3:4b; gemma3:12b; granite3.3:2b; granite3.3:8b; llama3.1:8b; llama3.2:1b; llama3.2:3b; mistral:7b; mistral-nemo:12b; mistral-small3.2:24b; PetrosStav/gemma3-tools:4b; phi3:3.8b; phi4:14b; gpt-3.5-turbo; gpt-4.1; gpt-4.1-mini; gpt-4o; gpt-4o-mini; claude-3-5-haiku-20241022; mistral-medium-2505; mistral-small-2506; mistral-large-2411
Reasoning-oriented LLMs	deepseek-r1:1.5b; deepseek-r1:7b; deepseek-r1:8b; deepseek-r1:14b; qwen3:4b; qwen3:8b; gpt-oss:20b; claude-sonnet-4-20250514; claude-opus-4-20250514; magistral-medium-2507; gpt-5-nano-2025-08-07; gpt-5-mini-2025-08-07; gpt-5-2025-08-07; claude-opus-4-1-20250805
Specialized models	led_large_16384_arxiv_summarization; medllama2:7b

With this selection, we covered models of different sizes and release periods, ensuring that both widely adopted systems and recent architectures were represented. Extraordinarily large models were not considered because their resource demands exceed what is practical for typical summarization pipelines and were beyond the resources available for this study.

These 50 diverse models were all tasked with generating summaries for each of the 1,000 abstracts in the dataset, resulting in 50,000 generated summaries available for evaluation.

2.3. Evaluation Metrics

As there is no single metric that can fully reflect summary quality, especially in the biomedical field where both coverage of key information and factual correctness are critical, we used a multitude of metrics grouped into three categories: traditional surface-level measures, embedding-based metrics, and performance-related measures that reflect the feasibility of using the methods in real-world applications. By combining all these metrics into one final overall score, we end up with a balanced benchmark value that reflects both summary quality and practical usability.

2.3.1. Surface-level Metrics

This group consists of metrics that compare the generated summaries with the reference summaries mainly at the word or phrase level. While they do not capture meaning beyond surface overlap, they remain common metrics in summarization research and provide a simple foundation for evaluation. We used three ROUGE variants (ROUGE-1, ROUGE-2, ROUGE-L) [25], BLEU [26], and METEOR [27]. ROUGE-1 and ROUGE-2

measure how many unigrams (single words) or bigrams (word pairs) from the reference appear in the generated output, while ROUGE-L identifies the longest sequence of words shared between the two. BLEU calculates how many n-grams in the output also occur in the reference, but it emphasizes precision rather than recall and applies a brevity penalty to counteract the tendency toward overly short summaries. METEOR extends n-gram matching by also considering word stems and synonyms, which makes it more tolerant to variations in wording. Together, these metrics offer a simple but transparent point of reference.

2.3.2. Embedding-based Metrics

To capture similarity beyond surface-level word overlap, we included a set of embedding-based metrics built on pre-trained transformer models. These methods generate vector representations of text, which allows them to capture similarity in meaning rather than just word overlap. We employed RoBERTa [28] and DeBERTa [29], two transformer-based models with strong performance across natural language processing tasks. In the context of summarization evaluation, they can be used to judge whether two summaries capture the same content even if phrased differently.

We also included all-mpnet-base-v2 [30], a transformer model fine-tuned for sentence similarity. Unlike RoBERTa and DeBERTa, which are primarily general-purpose encoders, MPNet was trained with a focus on aligning at the sentence-level. This focus makes it a useful complement to the other metrics, as it is particularly sensitive to whether the overall sense of a reference summary is preserved in the system output.

Finally, to evaluate factual consistency, we applied AlignScore [31], a metric designed to test whether the statements in a generated summary are supported by the source text. In contrast to the other metrics, we used AlignScore in a way where it does not compare the output to the reference summary but instead aligns it directly with the abstract, as factual accuracy can only be judged relative to the original input. This addition ensures that our evaluation is sensitive to errors and hallucinations that might otherwise be overlooked.

2.3.3. Performance Metrics

In addition to summary quality, we also considered practical aspects of model performance. Four measures were included: output token cost reflects the average length of generated summaries in tokens, as excessively long outputs increase runtime and resource requirements. Insufficient findings describe how often a model returned the predefined token 'INSUFFICIENT_FINDINGS' instead of producing a summary, capturing cases where it concluded the input did not contain substantive findings. Acceptance is the proportion of prompts for which a model produced an output, since some models occasionally failed to return a response. Finally, speed records the average time required to generate summaries, which is critical when processing large datasets.

These measures complement the quality metrics by addressing whether a method is not only accurate but also feasible to use in practice.

2.4. Benchmarking Framework

The benchmark was conducted using Python 3.12. Gold standard data were retrieved from open-access publications published by ScienceDirect and Cell Press through manual extraction of titles, abstracts, and highlight sections, along with metadata including publication URLs, identifiers, section types, and article types where available. All data were stored in machine-readable JSON format.

The framework was implemented using the Python standard library supplemented by several specialized packages: pandas [32] for data import and export, scikit-learn [33] for computing cosine similarities of embeddings and TF-IDF vectors, networkx [34] for

graph construction and PageRank algorithm [35]. Additional evaluation metrics were computed using NLTK [36] for METEOR and BLEU scores, ROUGE-score, BERT-score [37], AlignScore, and sentence-transformers [38] with the all-mpnet-base-v2 model.

Communication with proprietary closed-source LLMs was facilitated through the official Python APIs provided by Anthropic, Mistral AI, and OpenAI. Local LLM execution was performed on a workstation equipped with a NVIDIA RTX A4000 GPU (16GB VRAM) running Ollama as a backend service, accessed through its Python API along with the transformers library [39].

All LLMs were configured with a temperature parameter of 0.2 to optimize reproducibility while avoiding completely deterministic outputs. For the latest generation of OpenAI models featuring adaptive reasoning capabilities, the configuration was set to `text.verbosity = low` and `reasoning.effort = minimal`. The full set of parameters and prompts are documented in the `config.py` file in the repository.

2.5. Data Availability

The complete source code, documentation, gold standard dataset, and processed results are available at:

<https://www.github.com/Delta4AI/LLMTextSummarizationBenchmark>.

3. Results

This section may be divided by subheadings. It should provide a concise and precise description of the experimental results, their interpretation as well as the experimental conclusions that can be drawn.

3.1. Subsection

3.1.1. Subsubsection

Bulleted lists look like this:

- First bullet;
- Second bullet;
- Third bullet.

Numbered lists can be added as follows:

1. First item;
2. Second item;
3. Third item.

The text continues here.

3.2. Figures, Tables and Schemes

All figures and tables should be cited in the main text as Figure 1, Table 3, etc.



Figure 1. This is a figure. Schemes follow the same formatting.

Table 3. This is a table caption. Tables should be placed in the main text near to the first time they are cited.

Title 1	Title 2	Title 3
Entry 1	Data	Data
Entry 2	Data	Data ¹

¹ Tables may have a footer.

The text continues here (Figure 2 and Table 4).

167

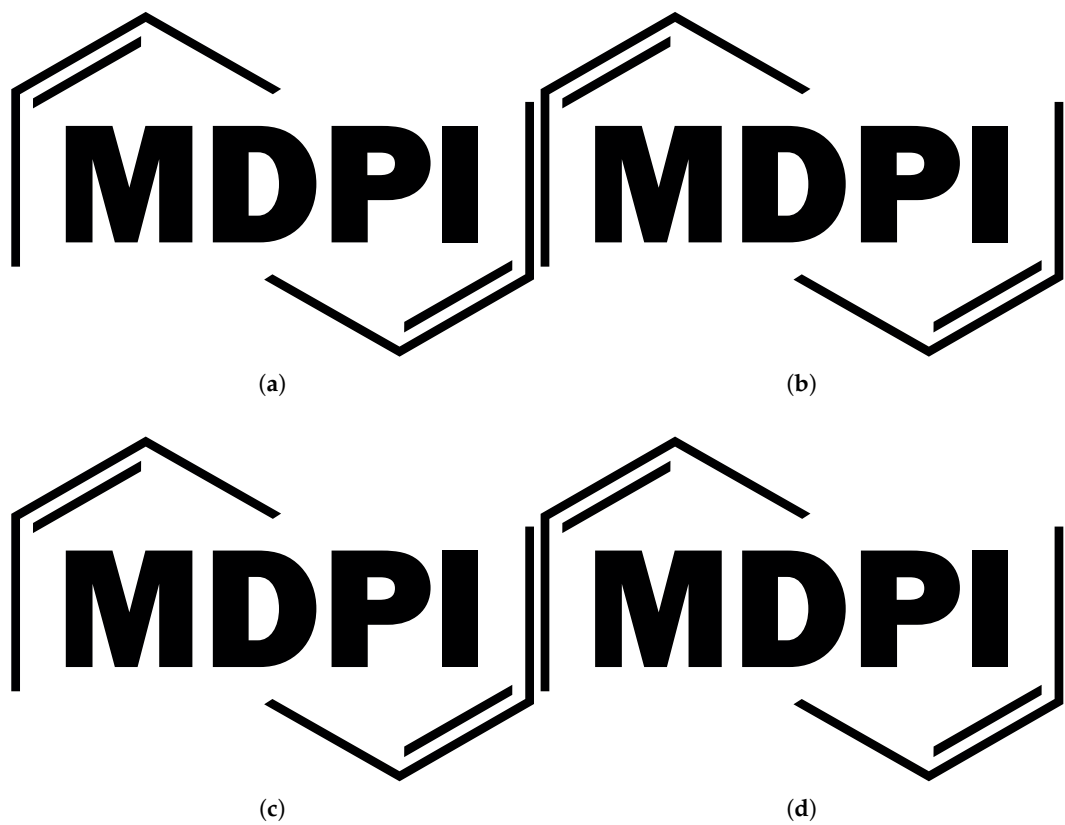


Figure 2. This is a wide figure. Schemes follow the same formatting. If there are multiple panels, they should be listed as: (a) Description of what is contained in the first panel. (b) Description of what is contained in the second panel. (c) Description of what is contained in the third panel. (d) Description of what is contained in the fourth panel. Figures should be placed in the main text near to the first time they are cited. A caption on a single line should be centered.

Table 4. This is a wide table.

Title 1	Title 2	Title 3	Title 4
Entry 1 *	Data	Data	Data
	Data	Data	Data
	Data	Data	Data
Entry 2	Data	Data	Data
	Data	Data	Data
	Data	Data	Data

* Tables may have a footer.

Text.
Text.

168

169

3.3. Formatting of Mathematical Components

This is the example 1 of equation:

$$a = 1, \quad (1)$$

the text following an equation need not be a new paragraph. Please punctuate equations as regular text.

This is the example 2 of equation:

$$a = b + c + d + e + f + g + h + i + j + k + l + m + n + o + p + q + r + s + t + u + v + w + x + y + z \quad (2)$$

Please punctuate equations as regular text. Theorem-type environments (including propositions, lemmas, corollaries etc.) can be formatted as follows:

Theorem 1. *Example text of a theorem.*

The text continues here. Proofs must be formatted as follows:

Proof of Theorem 1. Text of the proof. Note that the phrase “of Theorem 1” is optional if it is clear which theorem is being referred to. \square

The text continues here.

4. Discussion

Authors should discuss the results and how they can be interpreted from the perspective of previous studies and of the working hypotheses. The findings and their implications should be discussed in the broadest context possible. Future research directions may also be highlighted.

5. Conclusions

This section is not mandatory, but can be added to the manuscript if the discussion is unusually long or complex.

6. Patents

This section is not mandatory, but may be added if there are patents resulting from the work reported in this manuscript.

Author Contributions: For research articles with several authors, a short paragraph specifying their individual contributions must be provided. The following statements should be used “Conceptualization, X.X. and Y.Y.; methodology, X.X.; software, X.X.; validation, X.X., Y.Y. and Z.Z.; formal analysis, X.X.; investigation, X.X.; resources, X.X.; data curation, X.X.; writing—original draft preparation, X.X.; writing—review and editing, X.X.; visualization, X.X.; supervision, X.X.; project administration, X.X.; funding acquisition, Y.Y. All authors have read and agreed to the published version of the manuscript.”, please turn to the [CRediT taxonomy](#) for the term explanation. Authorship must be limited to those who have contributed substantially to the work reported.

Funding: Please add: “This research received no external funding” or “This research was funded by NAME OF FUNDER grant number XXX.” and “The APC was funded by XXX”. Check carefully that the details given are accurate and use the standard spelling of funding agency names at <https://search.crossref.org/funding>, any errors may affect your future funding.

Institutional Review Board Statement: In this section, you should add the Institutional Review Board Statement and approval number, if relevant to your study. You might choose to exclude this statement if the study did not require ethical approval. Please note that the Editorial Office might ask

you for further information. Please add “The study was conducted in accordance with the Declaration of Helsinki, and approved by the Institutional Review Board (or Ethics Committee) of NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving humans. OR “The animal study protocol was approved by the Institutional Review Board (or Ethics Committee) of NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving animals. OR “Ethical review and approval were waived for this study due to REASON (please provide a detailed justification).” OR “Not applicable” for studies not involving humans or animals.

Informed Consent Statement: Any research article describing a study involving humans should contain this statement. Please add “Informed consent was obtained from all subjects involved in the study.” OR “Patient consent was waived due to REASON (please provide a detailed justification).” OR “Not applicable” for studies not involving humans. You might also choose to exclude this statement if the study did not involve humans.

Written informed consent for publication must be obtained from participating patients who can be identified (including by the patients themselves). Please state “Written informed consent has been obtained from the patient(s) to publish this paper” if applicable.

Data Availability Statement: We encourage all authors of articles published in MDPI journals to share their research data. In this section, please provide details regarding where data supporting reported results can be found, including links to publicly archived datasets analyzed or generated during the study. Where no new data were created, or where data is unavailable due to privacy or ethical restrictions, a statement is still required. Suggested Data Availability Statements are available in section “MDPI Research Data Policies” at <https://www.mdpi.com/ethics>.

Acknowledgments: In this section you can acknowledge any support given which is not covered by the author contribution or funding sections. This may include administrative and technical support, or donations in kind (e.g., materials used for experiments). Where GenAI has been used for purposes such as generating text, data, or graphics, or for study design, data collection, analysis, or interpretation of data, please add “During the preparation of this manuscript/study, the author(s) used [tool name, version information] for the purposes of [description of use]. The authors have reviewed and edited the output and take full responsibility for the content of this publication.”

Conflicts of Interest: Declare conflicts of interest or state “The authors declare no conflicts of interest.” Authors must identify and declare any personal circumstances or interest that may be perceived as inappropriately influencing the representation or interpretation of reported research results. Any role of the funders in the design of the study; in the collection, analyses or interpretation of data; in the writing of the manuscript; or in the decision to publish the results must be declared in this section. If there is no role, please state “The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results”.

Abbreviations

The following abbreviations are used in this manuscript:

MDPI Multidisciplinary Digital Publishing Institute

DOAJ Directory of open access journals

TLA Three letter acronym

LD Linear dichroism

Appendix A

Appendix A.1

The appendix is an optional section that can contain details and data supplemental to the main text—for example, explanations of experimental details that would disrupt the flow of the main text but nonetheless remain crucial to understanding and reproducing the research shown; figures of replicates for experiments of which representative data are

shown in the main text can be added here if brief, or as Supplementary Data. Mathematical proofs of results not central to the paper can be added as an appendix.

Table A1. This is a table caption.

Title 1	Title 2	Title 3
Entry 1	Data	Data
Entry 2	Data	Data

Appendix B

All appendix sections must be cited in the main text. In the appendices, Figures, Tables, etc. should be labeled, starting with “A”—e.g., Figure A1, Figure A2, etc.

References

1. Nguyen, H.; Chen, H.; Pobbathi, L.; Ding, J. A Comparative Study of Quality Evaluation Methods for Text Summarization, 2024, [\[arXiv:cs/2407.00747\]](#). <https://doi.org/10.48550/arXiv.2407.00747>.
2. Elsevier. Highlights, 2024. Accessed: 2025-08-07.
3. Cell Press. Final Submission: Other Components: Highlights, 2024. Accessed: 2025-08-07.
4. Mihalcea, R.; Tarau, P. TextRank: Bringing Order into Text. In Proceedings of the Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing; Lin, D.; Wu, D., Eds., Barcelona, Spain, 2004; pp. 404–411.
5. Lewis, M.; Liu, Y.; Goyal, N.; Ghazvininejad, M.; Mohamed, A.; Levy, O.; Stoyanov, V.; Zettlemoyer, L. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. *CoRR* **2019**, *abs/1910.13461*, [\[1910.13461\]](#).
6. Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; Liu, P.J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* **2020**, *21*, 1–67.
7. Hasan, T.; Bhattacharjee, A.; Islam, M.S.; Mubasshir, K.; Li, Y.F.; Kang, Y.B.; Rahman, M.S.; Shahriyar, R. XL-Sum: Large-Scale Multilingual Abstractive Summarization for 44 Languages. In Proceedings of the Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, Online, 2021; pp. 4693–4703.
8. Zhang, J.; Zhao, Y.; Saleh, M.; Liu, P.J. PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization, 2019, [\[arXiv:cs.CL/1912.08777\]](#).
9. Team, G.; Kamath, A.; Ferret, J.; Pathak, S.; Vieillard, N.; Merhej, R.; Perrin, S.; Matejovicova, T.; Ramé, A.; Rivière, M.; et al. Gemma 3 Technical Report, 2025, [\[arXiv:cs.CL/2503.19786\]](#).
10. Mishra, M.; Stallone, M.; Zhang, G.; Shen, Y.; Prasad, A.; Soria, A.M.; Merler, M.; Selvam, P.; Surendran, S.; Singh, S.; et al. Granite Code Models: A Family of Open Foundation Models for Code Intelligence, 2024, [\[arXiv:cs.AI/2405.04324\]](#).
11. Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. The Llama 3 Herd of Models, 2024, [\[arXiv:cs.AI/2407.21783\]](#).
12. Jiang, A.Q.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chaplot, D.S.; de las Casas, D.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; et al. Mistral 7B, 2023, [\[arXiv:cs.CL/2310.06825\]](#).
13. Abdin, M.; Aneja, J.; Awadalla, H.; Awadallah, A.; Awan, A.A.; Bach, N.; Bahree, A.; Bakhtiari, A.; Bao, J.; Behl, H.; et al. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone, 2024, [\[arXiv:cs.CL/2404.14219\]](#).
14. Abdin, M.; Aneja, J.; Behl, H.; Bubeck, S.; Eldan, R.; Gunasekar, S.; Harrison, M.; Hewett, R.J.; Javaheripi, M.; Kauffmann, P.; et al. Phi-4 Technical Report, 2024, [\[arXiv:cs.CL/2412.08905\]](#).
15. Brown, T.B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. Language Models are Few-Shot Learners, 2020, [\[arXiv:cs.CL/2005.14165\]](#).
16. OpenAI; Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F.L.; Almeida, D.; Altenschmidt, J.; Altman, S.; et al. GPT-4 Technical Report, 2024, [\[arXiv:cs.CL/2303.08774\]](#).
17. Anthropic. Claude - Models overview, 2025. Accessed: 2025-09-24.
18. DeepSeek-AI; Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; et al. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, 2025, [\[arXiv:cs.CL/2501.12948\]](#).
19. Team, Q. Qwen3 Technical Report, 2025, [\[arXiv:cs.CL/2505.09388\]](#).
20. OpenAI. GPT-OSS: Open Source GPT Models, 2025. Accessed: 2025-09-23.
21. OpenAI. GPT-5 Models, 2025. Accessed: 2025-09-24.
22. Mistral-AI.; ; Rastogi, A.; Jiang, A.Q.; Lo, A.; Berrada, G.; Lample, G.; Rute, J.; Barmantlo, J.; Yadav, K.; et al. Magistral, 2025, [\[arXiv:cs.CL/2506.10910\]](#).

23. Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. Llama 2: Open Foundation and Fine-Tuned Chat Models, 2023, [[arXiv:cs.CL/2307.09288](https://arxiv.org/abs/2307.09288)].
24. Beltagy, I.; Peters, M.E.; Cohan, A. Longformer: The Long-Document Transformer. *arXiv:2004.05150* **2020**.
25. Lin, C.Y. ROUGE: A Package for Automatic Evaluation of Summaries. In Proceedings of the Text Summarization Branches Out, Barcelona, Spain, 2004; pp. 74–81.
26. Papineni, K.; Roukos, S.; Ward, T.; Zhu, W.J. Bleu: a Method for Automatic Evaluation of Machine Translation. In Proceedings of the Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics; Isabelle, P.; Charniak, E.; Lin, D., Eds., Philadelphia, Pennsylvania, USA, 2002; pp. 311–318. <https://doi.org/10.3115/1073083.1073135>.
27. Banerjee, S.; Lavie, A. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In Proceedings of the Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization; Goldstein, J.; Lavie, A.; Lin, C.Y.; Voss, C., Eds., Ann Arbor, Michigan, 2005; pp. 65–72.
28. Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; Stoyanov, V. RoBERTa: A Robustly Optimized BERT Pretraining Approach, 2019, [[arXiv:cs.CL/1907.11692](https://arxiv.org/abs/1907.11692)].
29. He, P.; Liu, X.; Gao, J.; Chen, W. DeBERTa: Decoding-enhanced BERT with Disentangled Attention, 2021, [[arXiv:cs.CL/2006.03654](https://arxiv.org/abs/2006.03654)].
30. Song, K.; Tan, X.; Qin, T.; Lu, J.; Liu, T.Y. MPNet: Masked and Permuted Pre-training for Language Understanding, 2020, [[arXiv:cs.CL/2004.09297](https://arxiv.org/abs/2004.09297)].
31. Zha, Y.; Yang, Y.; Li, R.; Hu, Z. AlignScore: Evaluating Factual Consistency with a Unified Alignment Function, 2023, [[arXiv:cs.CL/2305.16739](https://arxiv.org/abs/2305.16739)].
32. McKinney, W. Data Structures for Statistical Computing in Python. In Proceedings of the Proceedings of the 9th Python in Science Conference; van der Walt, S.; Millman, J., Eds., 2010, pp. 51–56.
33. Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **2011**, *12*, 2825–2830.
34. Hagberg, A.; Swart, P.; Schult, D. Exploring network structure, dynamics, and function using NetworkX. Technical report, Los Alamos National Lab.(LANL), Los Alamos, NM (United States), 2008.
35. Brin, S.; Page, L. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems* **1998**, *30*, 107–117. Proceedings of the Seventh International World Wide Web Conference, [https://doi.org/10.1016/S0169-7552\(98\)00110-X](https://doi.org/10.1016/S0169-7552(98)00110-X).
36. Bird, S.; Klein, E.; Loper, E. *Natural language processing with Python: analyzing text with the natural language toolkit*; "O'Reilly Media, Inc.", 2009.
37. Zhang, T.; Kishore, V.; Wu, F.; Weinberger, K.Q.; Artzi, Y. BERTScore: Evaluating Text Generation with BERT, 2020, [[arXiv:cs.CL/1904.09675](https://arxiv.org/abs/1904.09675)].
38. Reimers, N.; Gurevych, I. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, 2019, [[arXiv:cs.CL/1908.10084](https://arxiv.org/abs/1908.10084)].
39. Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; et al. Transformers: State-of-the-Art Natural Language Processing. In Proceedings of the Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, Online, 2020; pp. 38–45.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.